

# Performance analysis of electronic structure codes on HPC systems: A case study of SIESTA

Fabiano Corsetti<sup>1,\*</sup>

<sup>1</sup>*CIC nanoGUNE, 20018 Donostia-San Sebastián, Spain*

(Dated: April 23, 2014)

We report on scaling and timing tests of the SIESTA electronic structure code for *ab initio* molecular dynamics simulations using density-functional theory. The tests are performed on six large-scale supercomputers belonging to the PRACE Tier-0 network with four different architectures: Cray XE6, IBM BlueGene/Q, BullX, and IBM iDataPlex. We employ a systematic strategy for simultaneously testing weak and strong scaling, and propose a measure which is independent of the range of number of cores on which the tests are performed to quantify strong scaling efficiency as a function of simulation size. We find an increase in efficiency with simulation size for all machines, with a qualitatively different curve depending on the supercomputer topology, and discuss the connection of this functional form with weak scaling behaviour. We also analyze the absolute timings obtained in our tests, showing the range of system sizes and cores favourable for different machines. Our results can be employed as a guide both for running SIESTA on parallel architectures, and for executing similar scaling tests of other electronic structure codes.

## I. INTRODUCTION

The use of first principles atomistic simulations with density-functional theory<sup>1,2</sup> (DFT) has grown from a cottage industry in the early 1990s to a routine and integral part of many contemporary scientific disciplines, at the meeting point between condensed matter physics, physical chemistry, and the new range of nanosciences<sup>3,4</sup>. Potential practitioners have a large number of ready-made codes to choose from (see, e.g., Refs. [5–16]), which distinguish themselves in their licensing models, the range of features they offer, the specifics of the technical implementation, and, generally, where they lie on the (computational) cost–accuracy curve.

An important consideration for all modern DFT codes is their parallel scalability on high-performance computer (HPC) architectures, that open up the possibility of simulating very large physical systems entirely *ab initio*. Consequently, a substantial effort has gone into the development and optimization of many of these codes for the specific purpose of running on massively parallel systems<sup>11,17–28</sup>. Articles describing such developments typically illustrate the scaling performance of the code with an example of strong scaling<sup>9–27</sup>, i.e., the wall time speedup obtained for a simulation of fixed size over a range of number of cores. Less frequently, weak scaling performance (i.e., an increase of the problem size proportionally with the number of cores) is also shown<sup>22,25,28</sup>.

The use of a strong scaling example can be an effective way of giving a qualitative idea of the parallel efficiency of the code and the scale of problems which can realistically be solved with it. However, there are a number of issues in using the information as it is usually presented for extracting, even approximately, a generalized, quantitative measure of performance, such as could be used to attempt a comparison between codes.

Firstly, the range of cores over which this strong scaling is investigated is not fixed (as must be the case, since

time and memory requirements restrict the lower bound, and computational resources the upper bound). The significance of the demonstrated speedup depends crucially on the lower bound of this range; furthermore, the dependence is non-trivial. If we assume a constant rate of loss of efficiency as the parallelization is increased, a speedup of 3.9 when going from 8 to 32 cores should be better than a speedup of 3.8 when going from 2048 to 8192 cores; however, this is obviously not the case, as it is clear from experience that the actual rate increases significantly with the number of cores. Closer comparisons are even harder to judge: is a speedup of 3.8 between 512 and 2048 cores better or worse than a speedup of 3.7 between 1024 and 4096 cores? There is effectively no way to answer this question without making an assumption about how to model parallel performance in general. A well-known and popular, albeit extremely idealized, way to do is by Amdahl’s law<sup>29</sup>, that describes the overall speedup in terms of the parallelizable fraction of the code  $P$ . To the best of our knowledge, only one published strong scaling test for a DFT code<sup>20</sup> has reported on a fitted value for  $P$ .

Secondly, there is no standard physical system on which to test strong scaling. From the point of view of the material itself, this is somewhat understandable, as different codes specialize in different areas of modelling; a more fundamental problem, however, is that strong scaling efficiency changes with system size for a given material. Although some studies report system size dependent results<sup>9,12</sup>, this is generally not the case. How, then, to compare between, e.g., a strong scaling test on a 1532-atom carbon nanotube between 2048 and 32768 cores<sup>27</sup>, and one on a 1003-atom polyalanine peptide between 512 and 65536 cores<sup>23</sup>?

In this paper, we discuss these issues while reporting on tests of the parallel scaling performance of SIESTA<sup>7</sup>, a well-established DFT code based on norm-conserving pseudopotentials<sup>30</sup>, a basis of finite-range numerical atomic orbitals (NAOs), and an auxiliary real-space grid

System	Architecture	Topology	Processor type	Proc. speed (GHz)	Tot. cores	Tot. nodes	Cores/ node	Cores/ proc.	Mem./ core (GB)
Hermit	Cray XE6	3D torus	AMD Opteron	2.3	113664	3552	32	16	2/4
JUQUEEN	IBM BlueGene/Q	5D torus	IBM PowerPC A2	1.6	458752	28672	16	8	1
FERMI	IBM BlueGene/Q	5D torus	IBM PowerPC A2	1.6	163840	10240	16	8	1
Curie	BullX	Fat tree	Intel SandyBridge	2.7	80640	5040	16	8	4
SuperMUC	IBM iDataPlex	Fat tree	Intel SandyBridge	2.7	147456	9216	16	8	2
MareNostrum	IBM iDataPlex	Fat tree	Intel SandyBridge	2.6	48384	3024	16	8	2

TABLE I. Specifications for the six PRACE Tier-0 systems. We note that some systems include secondary types of nodes with different specifications; these are not listed here, and are not used for our tests.

for representing the electronic density. The tests are performed on six supercomputers (Table I), currently forming the network of Tier-0 systems of the Partnership for Advanced Computing in Europe<sup>31</sup> (PRACE). Our aims, therefore, are twofold:

- to give the most up-to-date, comprehensive and reliable results of the timing and scaling of SIESTA on modern HPC systems, so as to allow users of the code to calculate realistic timing estimates over a wide range of number of cores, and therefore plan how to make the best use of their computational resources;
- to propose a simple framework in which to analyze parallel scaling results for all electronic structure codes, arguing in particular for the use of Amdahl’s law to quantify strong scaling performance, and for the importance of investigating and reporting this measure as a function of system size.

## II. COMPUTATIONAL METHODS

Our scaling tests are performed on snapshots of liquid water in cubic boxes with periodic boundary conditions. This is the same system used previously for parallel benchmarking of the Quickstep<sup>9</sup> (CP2K) code; as noted by its authors, liquid water is ideal for this purpose, since boxes of any arbitrary number of molecules can be created while maintaining the same density and cell shape. Furthermore, the lack of crystalline symmetry and the 3D periodicity of the material ensure a sufficiently challenging task that we expect to give a fair idea of worst-case timings for most typical uses of the code, while the presence of a band gap ensures that we do not have to worry about convergence issues arising during the tests.

We simultaneously test weak and strong scaling (again similarly to Ref. [9]), by varying both the number of cores, from 32 to 4096 ( $N_c = 2^n, 5 \leq n \leq 12$ ), and the number of water molecules per core, from 1 to 32 ( $N_m/N_c = 2^n, 0 \leq n \leq 5$ ). The resulting suite of tests is shown in Fig. 1. The maximum system size tested is of 4096 water molecules (12288 atoms) for all values of  $N_m/N_c$ , except for one test of 8192 molecules (24576

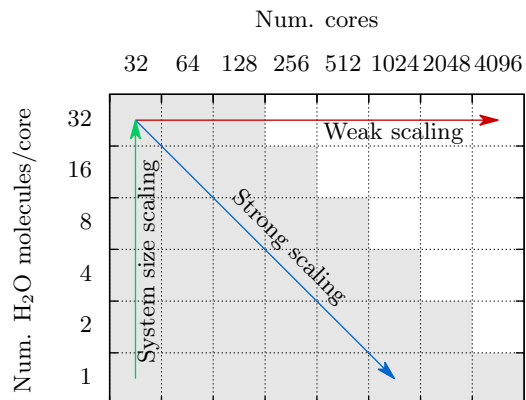


FIG. 1. Three types of scaling that can be investigated by systematically varying the number of molecules per core and the number of cores. The shaded cells show the suggested set of tests to perform on a typical HPC system.

atoms) on 8192 cores. We note, however, that due to the limited computational time available on each machine, not all tests are run on all machines. Weak scaling corresponds to moving perpendicular to the  $N_m/N_c$  axis, while strong scaling corresponds to moving diagonally. Instead, system size scaling (parallel to the  $N_m/N_c$  axis) does not explicitly test parallelization, although it is affected by it, as we shall discuss. The snapshots for all system sizes are extracted from classical molecular dynamics (MD) runs using the TIP4P force field<sup>32</sup> in the GROMACS<sup>33</sup> code, equilibrated to 300 K; the cell shape and volume are kept fixed at the experimental equilibrium density<sup>34</sup> (1.00 g/cm<sup>3</sup>).

The tests are performed at the  $\Gamma$  point only (multiple k points being almost embarrassingly parallel), using the semi-local PBE<sup>35</sup> functional for exchange and correlation (xc), a 150 Ry cutoff energy for the real-space auxiliary grid, and, unless otherwise stated, a double- $\zeta$  polarized basis<sup>36</sup> (d $\zeta$  + p), corresponding to 23 NAOs per water molecule; the fraction of occupied eigenstates is 4/23 ( $\sim 17\%$ ). All system sizes employ 13 self-consistent field (SCF) iterations to reach convergence.

We use the most recent development version of the code (`siesta-trunk-438`), available on the SIESTA website<sup>37</sup>. The tests are run with the code's default options for diagonalization, employing routines from the ScaLAPACK<sup>38</sup> library: the problem is first transformed from generalized to standard form by Cholesky factorization with the `pdpotrf` and `pdsygst` routines, and then the diagonalization itself is performed with the `pdsyevd` divide-and-conquer routine; finally, the back transform is performed with the `pdtrsm` routine. A 2D block-cyclic data distribution of the matrices is used, with the matrix dimension being an exact multiple of the block size in all cases (tests show the ideal block size to be equal to the number of orbitals per molecule).

We choose the standard solver for our tests, as this is currently by far the most widely used by the SIESTA community; however, we note that several new alternatives are being developed and tested: (i) a solver based on the orbital minimization method (OMM), which has already been demonstrated to exhibit better parallel scaling than explicit diagonalization up to 64 cores<sup>39</sup> (available in the development version of the code), (ii) two new solvers based on ScaLAPACK, the MRRR algorithm<sup>40</sup> and the ELPA library<sup>23</sup> (not yet released), and (iii) a solver based on the pole expansion and selected inversion method<sup>41</sup>, specifically designed for massively parallel architectures (not yet released). Finally, the original linear-scaling DFT method implemented in SIESTA is also in the process of being redesigned; in its current implementation it does not scale well on large clusters.

The code was compiled on each of the six machines listed in Table I using the native Fortran compiler and optimized linear algebra and communication libraries provided by the system administrators. The Intel compiler and MKL library are used for Intel-based machines (IBM iDataPlex and BullX architectures), the Cray compiler and ACML library for the AMD-based machine (Cray XE6 architecture), and the IBM XL compiler and ESSL library for the IBM PowerPC-based machines (IBM BlueGene/Q architecture). The MPI-2 libraries used are as follows: IBM MPI for SuperMUC, Open MPI for MareNostrum, BullX MPI for Curie, and MPICH2 for Hermit, JUQUEEN and FERMI.

### III. RESULTS AND DISCUSSION

#### A. Strong scaling

As previously mentioned, Amdahl's law provides a simple model of strong scaling. It states that

$$\mathbb{S}_1(N_c; S) = \frac{t_1}{t_{N_c}} = \frac{1}{S + \frac{1-S}{N_c}}, \quad (1)$$

where  $\mathbb{S}_1$  is the speedup obtained on  $N_c$  cores with respect to a serial run,  $t_1$  and  $t_{N_c}$  are the total execution times in serial and on  $N_c$  cores, respectively, and

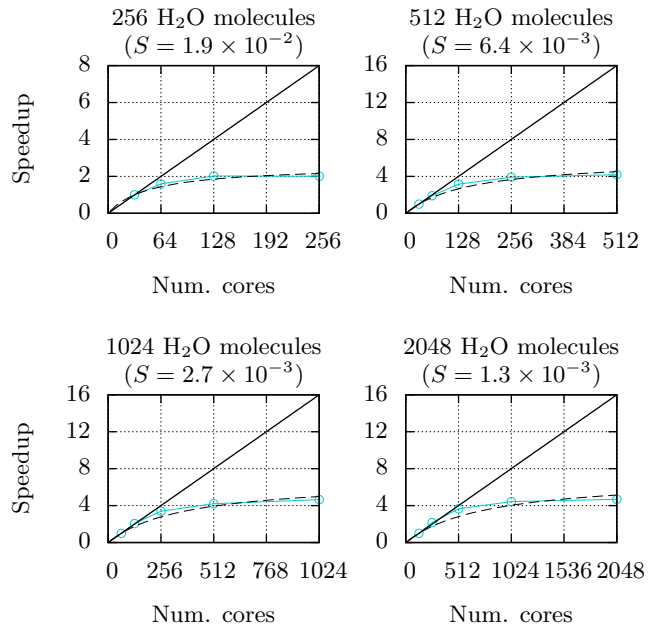


FIG. 2. Strong scaling on SuperMUC for four different system sizes. The full black lines gives the ideal scaling relative to the smallest system size. The fit to Amdahl's law is shown by the dashed black line, and the corresponding  $S$  value is given above the plot.

$S = 1 - P$  is the fraction of the code that is not parallelizable (we prefer using  $S$  instead of the more usual  $P$ , as the former tends to zero in the limit of ideal scaling). Since it is usually not possible in practice to measure  $t_1$  for large systems, it is useful to define the speedup with respect to a baseline number of cores  $b$  instead:

$$\mathbb{S}_b(N_c; S) = \frac{t_b}{t_{N_c}} = \left(\frac{N_c}{b}\right) \frac{S(b-1) + 1}{S(N_c-1) + 1}. \quad (2)$$

Using this equation, we can fit our strong scaling data over any arbitrary range of number of cores, and obtain a single value  $S$  that is in principle independent of this range, and which therefore defines the efficiency of the code for any value of  $N_c$  as  $1/(1 + S(N_c - 1))$  (see bottom panel, Fig. 3). The efficiency is invariantly 100% for a serial run, and decreases to zero as  $N_c \rightarrow \infty$ , since the execution time tends to a finite minimum value  $t_1 S$ .

It is important to note that the conventional interpretation for  $S$  and  $P$  is necessarily an over-simplification, and should not be taken too literally; nevertheless, Amdahl's law qualitatively reproduces some universal features of strong scaling, and is generally found to provide a good fit to real data. However, such a basic one-parameter model can only describe an average scaling trend, ignoring any system dependent effects that might favour particular values of  $N_c$ , e.g., differences in load balancing. Using a homogeneous, scalable system such as liquid water and a regular grid of tests as shown in Fig. 1 can be effective in minimizing these variations,

and therefore help to extract clearer general trends.

Using our timing tests for SIESTA on the six different machines, we can analyze the strong scaling of the code for system sizes ranging from 64 to 4096 water molecules; however, we restrict our fitting of  $S$  to systems with at least four data points ( $\geq 256$  molecules). As a representative example, Fig. 2 shows the speedup obtained on SuperMUC (IBM iDataPlex architecture) for four different system sizes, together with the curve fitted from Eq. 2. The resulting  $S$  values are robust to fitting over different ranges (within an order of magnitude), as is the trend of decreasing  $S$  with increasing system size. It is worth noting that this example clearly illustrates the difficulty in comparing between scaling tests using different ranges of number of cores: despite the steady increase in efficiency revealed by the  $S$  values, the speedups shown in the plots appear extremely similar due to the different baselines used.

The top panel of Fig. 3 summarizes the strong scaling results obtained for all six machines: the fitted value of  $S$  is given as a function of system size (i.e., the number of water molecules), for systems between 256 and 4096 molecules. Tests for smaller systems (32, 64, and 128 molecules) and larger ones (8192 molecules) are not represented, as there are insufficient data points for a reliable fit.

For all HPC systems,  $S$  is observed to decrease with the size of the physical system being simulated. This should not be surprising, as it is reasonable to expect efficiency to be related to the number of matrix elements/core (which in turn determines the ratio of intracore to inter-core operations), and, hence, that the larger the system being simulated, the larger the number of cores on which the calculation can be performed before the efficiency drops below a given threshold. However, the detailed form of this decrease depends on many factors related to the nature of the operations being performed and the computational architecture, and is therefore strongly dependent on the code and the HPC system used.

We can see some interesting distinctions in the  $S(N_m)$  curves for the six machines. There is a very close agreement between the three machines implementing torus topologies (Hermit, JUQUEEN, FERMI), despite Hermit being quite distinct from JUQUEEN and FERMI in most other respects, e.g., architecture type (Cray XE6 for the former, IBM BlueGene/Q for the latter) torus dimension, processor type and speed, number of cores per node and amount of memory per core. Instead, the three machines implementing fat tree topologies (Curie, SuperMUC, MareNostrum), even though they do not exhibit the same level of agreement amongst each other, give consistently higher  $S$  values than the torus machines.

Furthermore, despite the limited data available, our results suggest a qualitatively different form of the decrease of  $S$  with  $N_m$  for machines with torus and fat tree topologies. The former show an approximately linear decrease with slope  $B$  on the log-log scale ( $S \propto N_m^{-B}$ ), while the latter exhibit a slowing down of the rate of

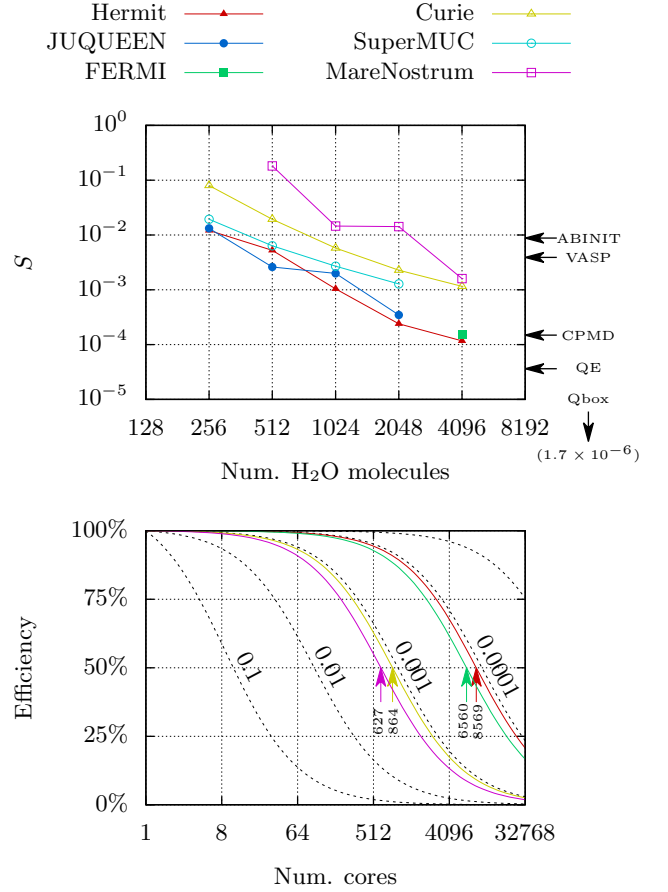


FIG. 3. Strong scaling and efficiency. Top panel:  $S$  value as a function of system size fitted to strong scaling data obtained with SIESTA on the six machines; also included are values calculated with other DFT codes for a single system size on IBM BlueGene architectures (ABINIT: 108 atoms, 1188 electrons, 3D system, 4 k points<sup>42</sup>; VASP: 87 atoms, 822 electrons, 2D system, 14 k points<sup>42</sup>; CPMD: 284 atoms, 1192 electrons, 3D system, k-point sampling unspecified<sup>42</sup>; QE: 1532 atoms, 5232 electrons, 1D system,  $\Gamma$  point<sup>27,42</sup>; Qbox: 1000 atoms, 12000 electrons, 3D system,  $\Gamma$  point<sup>11</sup>). Bottom panel: relationship between  $S$  and core hour efficiency as a function of the number of cores, for four different values of  $S$  given by the black dashed lines, and the fitted values of  $S$  obtained with SIESTA on four different machines for a system of 4096 water molecules; the number of cores at which the efficiency is equal to 50% is labelled in each case.

decrease. This is confirmed by fitting the data for each machine to a quadratic polynomial on the log-log scale; we find that the quadratic coefficient, positive in all cases, is an order of magnitude ( $\sim 3$ – $15$  times) smaller for the machines with torus topologies compared to those with fat tree topologies.

Using the  $S$  value as a measure of strong scaling, we can attempt a quantitative comparison between SIESTA and other DFT codes; this is given alongside our results in the top panel of Fig. 3. The fits are performed using



publicly available scaling test data for the codes, published on the website<sup>42</sup> of the FERMI IBM BlueGene/Q machine, which we also use for our tests of SIESTA; the same data for the Quantum ESPRESSO (QE) code has also been published in an article describing development work on the code<sup>27</sup>. The only exception is the Qbox code, for which we used previously published tests<sup>11</sup> performed on an IBM BlueGene/L machine. Where possible, we select tests performed  $\Gamma$ -only. All codes considered employ a plane-wave basis, in contrast to SIESTA's much smaller NAO basis.

It is important to stress that this comparison serves mainly to highlight the inadequacy of the available data; indeed, the change in  $S$  over more than three orders of magnitude for SIESTA at different system sizes is similar to the range spanned by the results obtained for the other codes, each available at a single system size. Both the system size and type vary greatly between codes, from an 87-atom 2D system for VASP to a 1532-atom 1D system for QE. Other important factors (k-point sampling, xc functional, basis accuracy, code optimization) are also not controlled for.

Nevertheless, Qbox stands out from all other codes for the impressive strong scaling performance demonstrated, with an  $S$  value more than an order of magnitude lower than that obtained by its closest competitor, QE, despite using a smaller system size (1000 atoms). Indeed, Qbox has been developed not only for massively parallel calculations in general, but specifically for running on IBM BlueGene architectures<sup>11,43</sup>; based on these results, it is the only DFT code to have demonstrated the potential to make efficient use ( $>50\%$ ) of the *entirety* of a large BlueGene machine such as JUQUEEN or FERMI for a single  $\Gamma$ -point calculation.

## B. System size scaling and absolute timings

Strong scaling is purely a test of parallel scalability, for which the code is, by definition, taken to be 100% efficient when run in serial. The results presented so far, therefore, contain no information about absolute timings. Although it is convenient to separate these two aspects of the code's performance, we should remember that the execution time is the *only* factor of importance to the end user. Therefore, strong scaling data on its own can sometimes be misleading, as a code that is very fast in serial but which exhibits poor strong scaling might nevertheless achieve a lower execution time on a medium-sized cluster than one that is very slow in serial but with exceptional scalability.

In order to extract a measure of absolute timing from our tests, we need to be able to effectively model system size scaling. For a conventional DFT code that calculates the eigenvalues and eigenvectors of the Kohn-Sham equation<sup>2</sup>, either by explicit diagonalization (as we do here) or by an iterative minimization algorithm, it is well known that the calculation time scales cubically with

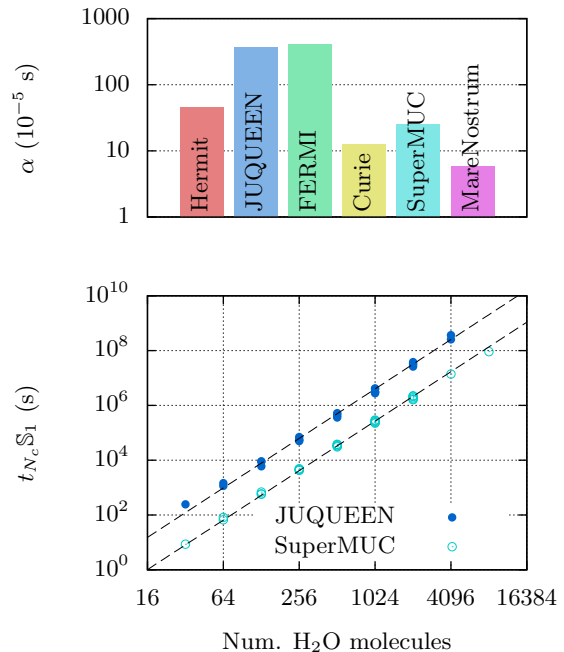


FIG. 4. Absolute timings on the six machines. Top panel: prefactor  $\alpha$  for the cubic scaling with system size of the execution time in serial for the self-consistent calculation of the liquid water system (13 SCF iterations). Bottom panel: two examples of the fitting of  $\alpha$  to absolute timing data, extrapolated for all number of cores to serial timings using Amdahl's law and a fitted analytical expression of the strong scaling performance as a function of system size.

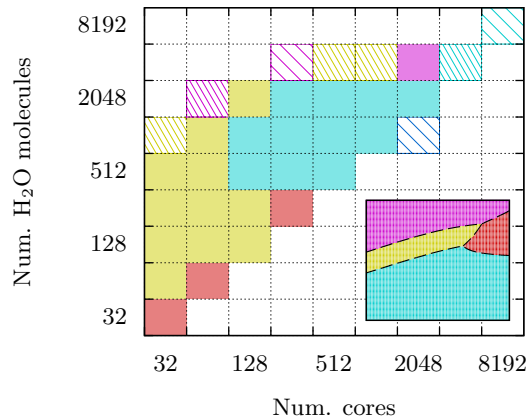


FIG. 5. Phase diagram of supercomputers. The machine with the lowest execution time is shown for a given system size and number of cores. The colours used are the same as those shown in the top panel of Fig. 4. Boxes with dashed lines indicate that the data for one or more machines is not available; sparse dashed lines indicate that only one machine was run with these parameters. The inset shows the idealized diagram over the same range, using the timing estimates given by the fits of  $S(N_m)$  and  $\alpha$ .

system size (i.e., the number of atoms/molecules/basis orbitals). Linear-scaling methods<sup>44</sup>, which make use of approximate spatial truncations based on the principle of electronic nearsightedness<sup>45</sup>, are also now well established and have been implemented in a number of popular codes.

The bottom panel of Fig. 4 shows the results of all timing tests performed on two machines, JUQUEEN and SuperMUC; we plot the execution time for all number of cores, extrapolated to that of a single core as  $t_{N_c} \mathbb{S}_1$  (from Eq. 1), against the system size (the number of water molecules  $N_m$ ). The estimated speedup  $\mathbb{S}_1$  is obtained for any value of  $N_m$  by using the fits of  $S(N_m)$  to the data in the top panel of Fig. 3, as described in the previous section. The resulting plot very clearly shows an almost pure cubic scaling with system size for both machines (the linear fits on the log-log scale have a fixed slope of 3). There is an excellent agreement in the extrapolated timings for each system size independently of the number of cores, and, even more encouragingly, our estimate of  $\mathbb{S}_1$  appears to be robust even when extrapolating beyond the range of  $N_m$  used in the fitting of  $S$ .

From these results, we can justify the use of a basic single-parameter model for system size scaling, of the form  $t_1 = \alpha N_m^3$ ; lower-order terms are negligible even for the smallest system sizes considered here; this is because all the default routines in SIESTA other than the diagonalization procedure itself are linear-scaling by design. Within an SCF iteration, the contribution from building the sparse Hamiltonian matrix only become comparable to diagonalization for very high values of the cutoff energy defining the real-space grid, or non-local xc functionals such as those including dispersion interactions. We note here that we have also analyzed the strong scaling of individual SIESTA modules, finding diagonalization to be the bottleneck within an SCF iteration, while Hamiltonian construction is very efficient when using the parallelization strategy for the grid operations of Sanz-Navarro *et al.*<sup>21</sup> (accessible via the flag `-DBSC_CELLXC` at compilation).

The parameter  $\alpha$ , obtained by the fits shown in the bottom panel of Fig. 4, can therefore be used to compare the speed of the various machines, independently of differences in scaling performance. The values of  $\alpha$  obtained for all six machines are shown in the top panel of Fig. 4. The large variation in  $\alpha$  over almost two orders of magnitude is a reflection not simply of the machines' processor speeds (listed in Table I), but also of numerous other interacting factors, such as the efficiency of the different compilers and libraries. In general, torus machines, which exhibit the best scaling, are predicted to be the slowest in serial, while fat tree machines, which do not scale as well in parallel, are predicted to be the fastest.

We can now calculate a rough estimate of the execution time on each machine for any number of water molecules on any number of cores, by using our fits of the function  $S(N_m)$  and the parameter  $\alpha$ , and, hence, build up a

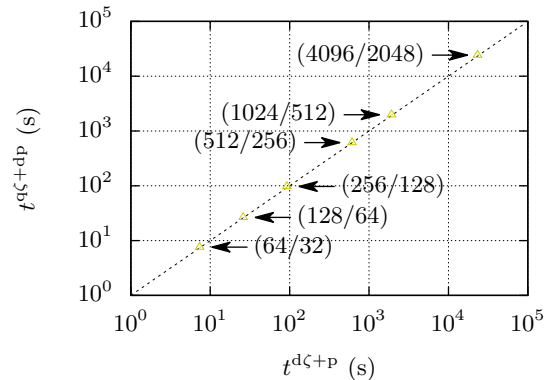


FIG. 6. Timing comparison on Curie for two different SIESTA basis sets. Each data point plots the execution time of a particular system size simulated with the  $d\zeta + p$  basis (23 NAOs/H<sub>2</sub>O molecule) against that of a *different* system size simulated with the  $q\zeta + dp$  basis (46 NAOs/H<sub>2</sub>O molecule), chosen so that the two systems have the same total number of basis orbitals. The two system sizes are shown in brackets ( $d\zeta + p/q\zeta + dp$ ); in each case, both simulations are performed on the same number of cores, equal to the number of molecules in the  $q\zeta + dp$  system.

$N_m$ - $N_c$  ‘phase diagram’ of the machine with the lowest execution time. This is shown in Fig. 5: the main panel compares real timings, while the inset uses the estimates based on our fits. The agreement is best for large system sizes and number of cores, with some discrepancies appearing for  $N_m \leq 256$ ; this is not surprising, due both to the extrapolation of  $S(N_m)$  to low values, and the fact that the timings are very close for more than one machine.

The machines which gives the lowest absolute timings over the entire tested range of  $N_c$  are overwhelmingly those with fat tree topologies, despite their inferior strong scaling performance with respect to torus machines. Two large regions can be clearly identified: Curie (BullX architecture) is the fastest machine for simulations with  $N_c \lesssim 128$ , while SuperMUC (IBM iDataPlex architecture) is the fastest above this value. There is some indication, confirmed by the model, that for large system sizes ( $N_m \gtrsim 4096$ ) MareNostrum (IBM iDataPlex architecture) becomes faster than both of these machines (this might seem surprising, since it has the lowest value of  $\alpha$ , and, hence, should be the fastest in serial at all system sizes; however, it also exhibits the worst parallel scaling, making it less efficient than other machines for parallel calculations on even very few cores at modest system sizes). It is only for extremely large  $N_c$  that the qualitatively different decrease in  $S(N_m)$  of the torus machines is predicted to lead to the lowest absolute timings, in particular for Hermit (Cray XE6 architecture), as it has a significantly lower  $\alpha$  value than the IBM BlueGene/Q machines.

Our fitted models for the six supercomputers can also

be used in a broader context, to estimate the execution time for any typical SIESTA calculation on HPC systems similar to the ones tested here. In fact, since the timing is dominated by the diagonalization procedure, especially for large system sizes, we can base our estimation on only two parameters, the total number of basis orbitals and the number of SCF iterations; we can safely neglect, to a first approximation, other parameters such as the number of electrons and the number and type of ions. This is illustrated in Fig. 6, in which we compare calculations using the standard  $d\zeta + p$  basis to ones using a larger  $q\zeta + dp$  basis<sup>46</sup> (twice the number of NAOs per water molecule) over a wide range of number of cores; as can be seen, timings on a given number of cores depend only on the total number of basis orbitals, and so a calculation using the larger basis takes the same time as one using the smaller basis with twice the system size. We note that this simple behaviour is due to the use of a solver which computes all eigenvalues by explicit diagonalization. Instead, solvers based on iterative minimization techniques (typically employed by plane-wave codes) scale only quadratically with the number of basis functions<sup>39</sup>; for such codes, we would expect the dependence of  $S(N_m)$  on basis size to be non-trivial. Unfortunately, we are not aware of published data for any other DFT code that could help in investigating this issue.

In order to allow SIESTA users to obtain absolute timing estimates for their parallel calculations, we have released a web applet<sup>47</sup> based on the model we have described and the quantitative data obtained from our scaling tests. We also include a version of the applet for offline use in the Supporting Information (`Code_S1.txt`); details of the fits and the final set of parameters for the six machines can be easily found within the code.

### C. Weak scaling

Finally, we briefly discuss the weak scaling behaviour demonstrated by the code. Weak scaling is of most interest to linear-scaling DFT codes<sup>22,25</sup>, for which the objective is to obtain a constant time-to-solution as the problem size is increased together with the number of cores (this is also known as Gustafson's law<sup>48</sup>). Cubic-scaling codes, instead, can achieve at best a quadratic weak scaling behaviour, which is rarely investigated<sup>28,39</sup>; nevertheless, it can provide useful information on the limits of efficiency of the code.

We find it convenient to plot the execution time divided by the square of the number of cores, so that ideal weak scaling behaviour will appear flat, analogously to the case of a linear-scaling code. A representative example for one machine, JUQUEEN, is shown in Fig. 7. Surprisingly, we observe better than ideal weak scaling, tending towards ideal as the number of cores is increased. The effect becomes more pronounced as the number of water molecules per core is decreased. These trends are almost perfectly reproduced by the timing estimates pro-

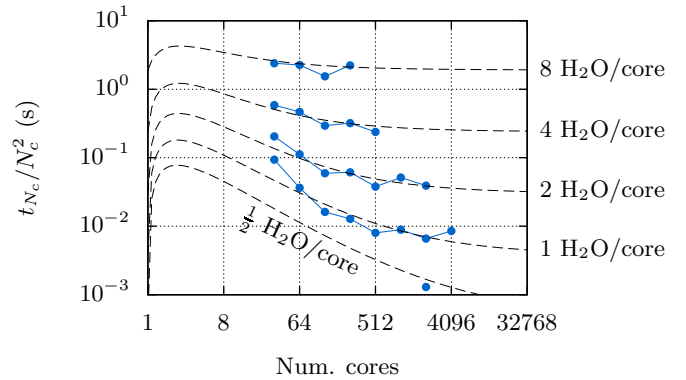


FIG. 7. Weak scaling on JUQUEEN for different numbers of water molecules per core. The execution time is divided by the square of the number of cores. The dashed lines show the estimates given by the fits of  $S(N_m)$  and  $\alpha$ .

vided by our combined modelling of strong scaling and system size scaling.

We can understand this behaviour as a change in efficiency (as defined in the bottom panel of Fig. 3) due to the interplay between the decrease of  $S$  with  $N_m$  and the increase of  $N_c$ . Similarly, it is interesting to note that system size scaling at a fixed number of cores  $> 1$  deviates from its ideal cubic behaviour in serial.

If we assume  $S(N_m)$  to be of the form  $AN_m^{-B}$ , it is easily verified from the model that the weak scaling behaviour will tend towards ideal for  $B \geq 1$ ; in the case of JUQUEEN, the fit gives a value of 1.8. This result applies equally to linear- and cubic-scaling codes, when using the appropriate definition of ideal weak scaling for each; indeed, a strikingly similar behaviour is reported for the linear-scaling Conquest code<sup>22</sup>. As noted previously, the three machines with fat tree topologies appear to exhibit a slowing down in the decrease of  $S$  with system size; although this should eventually make the weak scaling less than ideal, in practice it is not noticeable within the range of cores considered.

## IV. CONCLUSIONS

In this paper, we have investigated the performance of the SIESTA code on the six supercomputers of the PRACE Tier-0 network, currently amongst the largest in Europe. We propose a systematic investigation of parallel scaling using self-consistent calculations of snapshots of liquid water, varying both the number of cores on which the simulation is run and the number of water molecules per core; the largest simulation performed in our tests is of 8192 molecules on 8192 cores.

The results are analyzed using Amdahl's law to fit the data for each system size, providing a quantitative estimate of the code's efficiency over all number of cores

based on a single parameter  $S$ ; the scaling performance of the code, therefore, is completely described by the curve of  $S$  as a function of system size. We find a qualitative difference in this curve depending on the topology of the connections between nodes in the supercomputer, with machines implementing torus topologies demonstrating a better scalability to large system sizes than those implementing fat tree topologies. Despite this, however, the latter group is shown to give lower absolute timings for almost all simulations within the tested range, as the performance on individual cores is significantly faster; furthermore, such architectures tend to offer a larger amount of memory per core, which can become an important issue either when running on few cores, or as the size of the simulation is increased (the memory requirements scale approximately quadratically with system size).

Combining Amdahl's law for strong scaling with a basic one-parameter model for system size scaling, both of which are fitted to the data provided by our tests, we can calculate a simple estimate of the execution time on a given number of cores for a generic total energy calculation with SIESTA; a new web applet<sup>47</sup> developed in conjunction with the paper allows users of the code to employ this model for planning their projects on parallel architectures. An estimate of the memory requirement per core is also included.

Throughout the paper we have emphasized potential points of comparison with other DFT and electronic structure codes. Investigating and reporting  $S(N_m)$  curves for different HPC systems could provide valuable information to practitioners in the field, as well as for the ongoing development of the codes themselves. Care

must be taken, however, when interpreting the results of comparisons based on strong scaling data, due to the fundamental differences between codes. Basis sets offer perhaps the most important example: is it meaningful to compare the strong scaling performance of a localized-orbital code and a plane-wave code for the same physical system? It is clear that  $S$  varies with basis size, and so is crucially dependent in both cases on the precision level of the calculation; even disregarding the technical challenges involved<sup>46</sup>, attempting to equate the two bases is not necessarily appropriate, as the codes are designed from the outset to be used with different aims. For this reason, we suggest that the best approach should not be overly competitive; rather, the objective should be to report on calculations using the typical setup appropriate for each code (e.g., the default d $\zeta$  + p basis for SIESTA), or possibly a range of different setups, as this will provide the most useful information for its users.

## ACKNOWLEDGMENTS

We thank Emilio Artacho, Alberto Garcia, and Georg Huhs for useful discussions. We acknowledge PRACE for awarding us access to the following resources: Hermit based in Germany at the High Performance Computing Center Stuttgart (HLRS), JUQUEEN based in Germany at the Jülich Supercomputing Centre, FERMI based in Italy at the CINECA SuperComputing Application and Innovation Department, Curie based in France at the Très Grand Centre de Calcul du CEA (TGCC), SuperMUC based in Germany at the Leibniz Supercomputing Centre, and MareNostrum based in Spain at the Barcelona Supercomputing Center.

- 
- \* E-mail: f.corsetti@nanogune.eu
- <sup>1</sup> P. Hohenberg and W. Kohn, Phys. Rev. **136**, B864 (1964).
  - <sup>2</sup> W. Kohn and L. J. Sham, Phys. Rev. **140**, A1133 (1965).
  - <sup>3</sup> J. Hafner, C. Wolverton, and G. Ceder, MRS Bull. **31**, 659 (2006).
  - <sup>4</sup> N. Marzari, MRS Bull. **31**, 681 (2006).
  - <sup>5</sup> G. Kresse and J. Furthmüller, Phys. Rev. B **54**, 11169 (1996).
  - <sup>6</sup> B. Delley, J. Chem. Phys. **113**, 7756 (2000).
  - <sup>7</sup> J. M. Soler, E. Artacho, J. D. Gale, A. García, J. Junquera, P. Ordejón, and D. Sánchez-Portal, J. Phys.: Condens. Matter **14**, 2745 (2002).
  - <sup>8</sup> S. J. Clark, M. D. Segall, C. J. Pickard, P. J. Hasnip, M. I. J. Probert, K. Refson, and M. C. Payne, Z. Kristallogr. **220**, 567 (2005).
  - <sup>9</sup> J. VandeVondele, M. Krack, F. Mohamed, M. Parrinello, T. Chassaing, and J. Hutter, Comput. Phys. Commun. **167**, 103 (2005).
  - <sup>10</sup> C.-K. Skylaris, P. D. Haynes, A. A. Mostofi, and M. C. Payne, J. Chem. Phys. **122**, 084119 (2005).
  - <sup>11</sup> F. Gygi, IBM J. Res. Dev. **52**, 1 (2008).
  - <sup>12</sup> L. Genovese, A. Neelov, S. Goedecker, T. Deutsch, S. A. Ghasemi, A. Willand, D. Caliste, O. Zilberberg, M. Rayson, A. Bergman, and R. Schneider, J. Chem. Phys. **129**, 014109 (2008).
  - <sup>13</sup> P. Giannozzi, S. Baroni, N. Bonini, M. Calandra, R. Car, C. Cavazzoni, D. Ceresoli, G. L. Chiarotti, M. Cococcioni, I. Dabo, A. Dal Corso, S. de Gironcoli, S. Fabris, G. Fratesi, R. Gebauer, U. Gerstmann, C. Gougousis, A. Kokalj, M. Lazzeri, L. Martin-Samos, N. Marzari, F. Mauri, R. Mazzarello, S. Paolini, A. Pasquarello, L. Paulatto, C. Sbraccia, S. Scandolo, G. Sciauzero, A. P. Seitsonen, A. Smogunov, P. Umari, and R. M. Wentzcovitch, J. Phys.: Condens. Matter **21**, 395502 (2009).
  - <sup>14</sup> V. Blum, R. Gehrke, F. Hanke, P. Havu, V. Havu, X. Ren, K. Reuter, and M. Scheffler, Comput. Phys. Commun. **180**, 2175 (2009).
  - <sup>15</sup> X. Gonze, B. Amadon, P.-M. Anglade, J.-M. Beuken, F. Bottin, P. Boulanger, F. Bruneval, D. Caliste, R. Caracas, M. Côté, T. Deutsch, L. Genovese, P. Ghosez, M. Giantomassi, S. Goedecker, D. R. Hamann, P. Hermet, F. Jollet, G. Jomard, S. Leroux, M. Mancini, S. Mazevet, M. J. T. Oliveira, G. Onida, Y. Pouillon, T. Rangel, G.-M. Rignanese, D. Sangalli, R. Shaltaf, M. Torrent, M. J. Ver-



- straete, G. Zerah, and J. W. Zwanziger, *Comput. Phys. Commun.* **180**, 2582 (2009).
- <sup>16</sup> J. Enkovaara, C. Rostgaard, J. J. Mortensen, J. Chen, M. Dulak, L. Ferrighi, J. Gavnholt, C. Glinsvad, V. Haikola, H. A. Hansen, H. H. Kristoffersen, M. Kuisma, A. H. Larsen, L. Lehtovaara, M. Ljungberg, O. Lopez-Acevedo, P. G. Moses, J. Ojanen, T. Olsen, V. Petzold, N. A. Romero, J. Stausholm-Møller, M. Strange, G. A. Tritsarlis, M. Vanin, M. Walter, B. Hammer, H. Häkkinen, G. K. H. Madsen, R. M. Nieminen, J. K. Nørskov, M. Puska, T. T. Rantala, J. Schiøtz, K. S. Thygesen, and K. W. Jacobsen, *J. Phys.: Condens. Matter* **22**, 253202 (2010).
  - <sup>17</sup> M. Plummer, J. Hein, M. F. Guest, K. J. D'Mellow, I. J. Bush, K. Refson, G. J. Pringle, L. Smith, and A. Trew, *J. Mater. Chem.* **16**, 1885 (2006).
  - <sup>18</sup> F. Bottin, S. Leroux, A. Knyazev, and G. Zerah, *Comp. Mater. Sci.* **42**, 329 (2008).
  - <sup>19</sup> L. Genovese, M. Ospici, T. Deutsch, J.-F. Méhaut, A. Neelov, and S. Goedecker, *J. Chem. Phys.* **131**, 034103 (2009).
  - <sup>20</sup> N. D. M. Hine, P. D. Haynes, A. A. Mostofi, C.-K. Skylaris, and M. C. Payne, *Comput. Phys. Commun.* **180**, 1041 (2009).
  - <sup>21</sup> C. F. Sanz-Navarro, R. Grima, A. García, E. A. Bea, A. Soba, J. M. Cela, and P. Ordejón, *Theor. Chem. Acc.* **128**, 825 (2010).
  - <sup>22</sup> D. R. Bowler and T. Miyazaki, *J. Phys.: Condens. Matter* **22**, 074207 (2010).
  - <sup>23</sup> T. Auckenthaler, V. Blum, H.-J. Bungartz, T. Huckle, R. Johanni, L. Krmer, B. Lang, H. Lederer, and P. Willems, *Parallel Comput.* **37**, 783 (2011).
  - <sup>24</sup> A. Maniopoulou, E. R. Davidson, R. Grau-Crespo, A. Walsh, I. J. Bush, C. R. A. Catlow, and S. M. Woodley, *Comput. Phys. Commun.* **183**, 1696 (2012).
  - <sup>25</sup> J. VandeVondele, U. Borštnik, and J. Hutter, *J. Chem. Theory Comput.* **8**, 3565 (2012).
  - <sup>26</sup> M. Hacene, A. Anciaux-Sedrakian, X. Rozanska, D. Klahr, T. Guignon, and P. Fleurat-Lessard, *J. Comput. Chem.* **33**, 2581 (2012).
  - <sup>27</sup> N. Varini, D. Ceresoli, L. Martin-Samos, I. Girotto, and C. Cavazzoni, *Comput. Phys. Commun.* **184**, 1827 (2013).
  - <sup>28</sup> S. Hakala, V. Havu, J. Enkovaara, and R. Nieminen, in *Applied parallel and scientific computing, Lecture notes in computer science*, Vol. 7782, edited by P. Manninen and P. Öster (Springer, Heidelberg, 2013) pp. 63–76.
  - <sup>29</sup> G. M. Amdahl, in *Proceedings of the April 18–20, 1967, spring joint computer conference*, AFIPS '67 (Spring), Vol. 30 (ACM, New York, 1967) pp. 483–485.
  - <sup>30</sup> N. Troullier and J. L. Martins, *Phys. Rev. B* **43**, 1993 (1991).
  - <sup>31</sup> URL: <http://www.prace-ri.eu/>.
  - <sup>32</sup> W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, *J. Chem. Phys.* **79**, 926 (1983).
  - <sup>33</sup> H. J. C. Berendsen, D. van der Spoel, and R. van Drunen, *Comput. Phys. Commun.* **91**, 43 (1995).
  - <sup>34</sup> W. Wagner and A. Pruß, *J. Phys. Chem. Ref. Data* **31**, 387 (2002).
  - <sup>35</sup> J. P. Perdew, K. Burke, and M. Ernzerhof, *Phys. Rev. Lett.* **77**, 3865 (1996).
  - <sup>36</sup> J. Wang, G. Román-Pérez, J. M. Soler, E. Artacho, and M.-V. Fernández-Serra, *J. Chem. Phys.* **134**, 024516 (2011).
  - <sup>37</sup> URL: <http://departments.icmab.es/leem/siesta/>.
  - <sup>38</sup> L. S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley, *ScaLAPACK users' guide* (Society for Industrial and Applied Mathematics, Philadelphia, 1997).
  - <sup>39</sup> F. Corsetti, *Comput. Phys. Commun.* **185**, 873 (2014).
  - <sup>40</sup> D. Antonelli and C. Voemel, *LAPACK working note 168: pdsyevr. ScaLAPACK's parallel MRRR algorithm for the symmetric eigenvalue*, Tech. Rep. UCB/CSD-05-1399 (EECS Department, University of California, Berkeley, 2005).
  - <sup>41</sup> L. Lin, M. Chen, C. Yang, and L. He, *J. Phys.: Condens. Matter* **25**, 295501 (2013).
  - <sup>42</sup> URL: <http://www.hpc.cineca.it/content/fermi-software-benchmarks/>.
  - <sup>43</sup> F. Gygi, E. W. Draeger, M. Schulz, B. R. de Supinski, J. A. Gunnels, V. Austel, J. C. Sexton, F. Franchetti, S. Kral, C. W. Ueberhuber, and J. Lorenz, in *Proceedings of the 2006 ACM/IEEE conference on supercomputing*, SC '06 (ACM, New York, 2006) pp. 45–52.
  - <sup>44</sup> D. R. Bowler and T. Miyazaki, *Rep. Prog. Phys.* **75**, 036503 (2012).
  - <sup>45</sup> W. Kohn, *Phys. Rev. Lett.* **76**, 3168 (1996).
  - <sup>46</sup> F. Corsetti, M.-V. Fernández-Serra, J. M. Soler, and E. Artacho, *J. Phys.: Condens. Matter* **25**, 435504 (2013).
  - <sup>47</sup> URL: <http://departments.icmab.es/leem/siesta/siestimator/>.
  - <sup>48</sup> J. L. Gustafson, *Commun. ACM* **31**, 532 (1988).